



# Energy-saving Pushing Based on Personal Interest and Context Information

Chuting Yao, Binqiang Chen and Chenyang Yang

Beihang University, Beijing, China  
Email: {ctyao, chenbq, cyyang}@buaa.edu.cn

Gang Wang

NEC Labs, China  
Email: wang\_gang@nec.cn

**Abstract**—Pushing files to users based on predicting the personal interest of each user may provide higher throughput gain than broadcasting popular files to users based on their common interests. However, the energy consumed at base station for pushing files individually to each user is also higher than broadcast. In this paper, we propose an energy-saving transmission strategy for pre-downloading the files to each user by exploiting the excess resources in the network during off-peak time. Specifically, a power allocation and scheduling algorithm is designed aimed to minimize the extra energy consumed for pushing, where network and user level context information are exploited. Simulation results show that when the energy of both content placement and content delivery is taken into account, the proposed unicast strategy consumes less energy and achieves higher throughput than broadcasting when the files popularity is not uniform and the personal interest prediction is with less uncertainty.

**Index Terms**—Pushing, unicast and broadcast, energy-saving

## I. INTRODUCTION

Since a large amount of traffic is generated by a few popular contents, proactively caching contents at base stations (BSs) or even at users can reduce backhaul cost [1] and end-to-end delay [2], improve energy efficiency [3–5], and boost network throughput by offloading [1, 6, 7], which is a promising way to support the explosively increasing traffic demands of fifth generation (5G) cellular networks with low cost.

Proactive caching consists of *content placement* before the actual demand arrives and *content delivery* after the user requests the files. Placing contents to users according to the predicted demands can offload wireless traffic in peak time to off-peak time. If a user can find its requested contents in local cache, it is no need for a BS to deliver the contents to the cache hit user. Otherwise, the content delivery consumes the spectrum and the energy resources at the BS.

With the predicted personal interest of each user, pushing files to an individual user has long been regarded as a technique for improving the user experience [8], where the content placement is accomplished by *unicast* when the channel of the user is in good condition [8, 9]. With the predicted common interest of the users in a cell (i.e., file popularity), the same set of popular contents can be pre-cached at the users by *broadcast* [6, 7]. Because the demand statistics of the users in a cell differ from that of each user

[2], pushing the possibly interested files to a user during the off-peak time before its actual demand arrives may provide higher throughput and lower energy in content delivery than broadcasting popular contents to all users thanks to the higher hit rate. However, because unicast has to be applied due to the different interests of multiple users, the content placement by pushing may consume more energy than broadcast. To enjoy the high throughput gain of pushing without incurring high cost, which can be captured by the energy consumed at the BS for both content placement and content delivery, it is important to develop energy-saving transmission strategy for pre-downloading the files to each user.

In this paper, we strive to optimize power allocation and user scheduling that minimizes the energy consumed by the BSs for the unicast content placement. In order to exploit the excess network resource, we develop an algorithm based on the user and network level context information. Simulation results demonstrate that the proposed strategy provides high network throughput and low energy consumption than broadcasting popular files to all users in a cell, where the energy consumed by both content placement and content delivery is taken into account, especially when the average arrival rate of the content delivery request is high.

## II. SYSTEM MODEL

Consider  $M$  small cells hexagonally deployed in a macro cell, where each small BS (SBS) is equipped with  $N_t$  antennas. The maximal transmit power of each SBS is  $p_{\max}$ , and the maximal bandwidth is  $W_{\max}$ .

### A. Traffic Model

Consider that the SBSs serve  $K$  randomly arrived MSs requesting files of  $B$  bits in total from a catalog with a time-slotted fashion, where the duration of each time slot is  $\Delta_t$ . At the same time, a given portion of the resources are reserved for each real-time (RT) request such as phone call, which has high priority. The content delivery can only use the residual resources at each SBS, which is time-varying. Denote the transmit power and bandwidth occupied by the RT traffic of the  $i$ th SBS in the  $t$ th time slot as  $p_{i,RT}^t$  and  $W_{i,RT}^t$ .

The arrival of the content delivery requests is often bursty owing to the human behavior, which exhibits peaks (e.g., during the lunch and dinner times). To capture this feature,

we assume that the requests of content delivery only arrive at the peak time. For mathematical tractability, assume that only one mobile station (MS) can be served by a SBS in each time slot and the MS only accesses to the closest SBS. Denote  $s_{i,k}^t \in \{1, 0\}$ ,  $i = 1, \dots, M$ ,  $k = 1, \dots, K$  as the scheduling indicator. When the  $i$ th SBS (denoted as BS <sub>$i$</sub> ) schedules the  $k$ th user (denoted as MS <sub>$k$</sub> ) in the  $t$ th time slot,  $s_{i,k}^t = 1$ , otherwise  $s_{i,k}^t = 0$ . Then,  $\sum_{i=1}^M s_{i,k}^t \leq 1$  and  $\sum_{k=1}^K s_{i,k}^t \leq 1$ .

Assume that the macro BS (MBS) can predict the interest of each MS, and gather the user and network level context information [10] from each SBS and each MS. With these information, the MBS can make the resource usage plan for each SBS to pre-download (also refer to as push in the sequel) the possibly interested contents to each MS during the off-peak time before the MS's request arrives.

### B. Channel Model

We divide the off-peak time into  $T_f$  frames. Each frame is further divided into  $T_s$  time slots. Hence, the entire off-peak time contains  $T \triangleq T_f T_s$  time slots. Due to the user mobility, the large-scale fading gain may vary among different frames. The small-scale fading is modelled as block fading, which may vary among time slots and remains constant in each time slot.

The received signal of MS <sub>$k$</sub>  in the  $t$ th time slot is

$$y_k^t = \sum_{i=1}^M s_{i,k}^t \sqrt{\alpha_k^{\lceil \frac{T}{T_s} \rceil}} (\mathbf{h}_k^t)^H \mathbf{w}_k^t \sqrt{p_k^t} x_k^t + n_k^t, \quad (1)$$

where  $x_k^t$  is the transmit symbol with  $\mathbb{E}\{|x_k^t|^2\} = 1$ ,  $p_k^t$  is the transmit power,  $\mathbf{w}_k^t \in \mathbb{C}^{N_t \times 1}$  is the beamforming vector,  $\mathbf{h}_k^t \in \mathbb{C}^{N_t \times 1}$  is the independent and identically distributed (i.i.d.) Rayleigh fading channel vector,  $\alpha_k^{\lceil \frac{T}{T_s} \rceil}$  is the corresponding large-scale fading gain between the user and the closest SBS, and  $n_k^t$  is the noise with variance  $\sigma^2$ .  $\mathbb{E}\{\cdot\}$  represents expectation, and  $\lceil \cdot \rceil$  is the ceiling function. Given that MS <sub>$k$</sub>  is scheduled only by one SBS in each time slot, maximum ratio transmission is optimal, i.e.,  $\mathbf{w}_k^t = \mathbf{h}_k^t / \|\mathbf{h}_k^t\|$ , where  $\|\cdot\|$  denotes Euclidean norm.

In the  $t$ th time slot, the achievable rate of MS <sub>$k$</sub>  in nats is

$$R_k^t = \sum_{i=1}^M s_{i,k}^t W_i^t \ln(1 + g_k^t p_k^t), \quad (2)$$

where  $W_i^t \triangleq W_{\max} - W_{i,\text{RT}}^t$  is the residual bandwidth available for pushing files at BS <sub>$i$</sub>  in the  $t$ th time slot,  $g_k^t \triangleq \alpha_k^{\lceil \frac{T}{T_s} \rceil} \|\mathbf{h}_k^t\|^2 / \sigma^2 = \alpha_k^{\lceil \frac{T}{T_s} \rceil} \|\mathbf{h}_k^t\|^2 / (N_0 W_i^t)$  is the equivalent channel gain and  $N_0$  is noise power spectrum density.

### C. Power Model

Assume that a SBS can be switched into sleep mode when the SBS has no traffic to serve. The pre-downloading power consumed by the SBSs at the  $t$ th time slot contains the transmit power for multiple users and the *extra* circuit power consumed for operating the SBS to push the files, which can be modeled as [11],

$$p_C^t = \sum_{k=1}^K \frac{1}{\xi} p_k^t + \sum_{i=1}^M \mathbf{1}(p_{i,\text{RT}}^t = 0) \mathbf{1}(\sum_{k=1}^K s_{i,k}^t > 0) (p_{\text{act}} - p_{\text{sle}}), \quad (3)$$

where  $\xi$  is the power amplifier efficiency,  $p_{\text{act}}$  and  $p_{\text{sle}}$  are the circuit power consumptions when the SBS is in active and sleep modes, respectively, and  $\mathbf{1}(x) = 1$  when the event  $x$  is true, otherwise,  $\mathbf{1}(x) = 0$ .

## III. ENERGY-SAVING PRE-DOWNLOADING

Based on the predicted personal interests of each user, the SBSs can push the contents to each user during the off-peak time to exploit the excess spectrum resource. However, energy will be the wasted if the pre-downloaded files are not requested by the users. To circumvent this problem, we propose a context-aware transmission strategy for pre-downloading that minimizes the extra energy consumed for pushing. For easy elaboration, we start by single-user case, and then extend the strategy to multi-user case.

### A. Single-user Pre-downloading

We first optimize the transmission strategy to pre-download the  $B$  bits contents to one user (say MS <sub>$k$</sub> ) in the off-peak time with duration  $T$ . Specifically, we optimize  $s_{i,k}^t$  and  $p_k^t$  in each time slot. For notational simplicity, we omit the subscripts  $i$  and  $k$  in this subsection. Since the user is only accessed to its closest SBS in each time slot,  $s^t = 1$  when  $p^t > 0$  and  $s^t = 0$  otherwise. Hence, we only need to optimize  $p^t$ .

If all the information of  $g^t, W_{\text{RT}}^t, p_{\text{RT}}^t, t = 1, \dots, T$  is available at the MBS in off-peak time, the MBS can optimize the resource allocation for the SBS to minimize the pre-downloading energy consumption from the following problem,

$$\min_{p^1, \dots, p^T} \frac{1}{T} \sum_{t=1}^T \left( \frac{1}{\xi} p^t + \mathbf{1}(p_{\text{RT}}^t = 0) \mathbf{1}(p^t > 0) (p_{\text{act}} - p_{\text{sle}}) \right) \quad (4a)$$

$$\text{s.t. } \frac{1}{T} \sum_{t=1}^T W^t \ln(1 + g^t p^t) = \frac{B \ln 2}{T \Delta t}, \quad (4b)$$

$$p^t \geq 0, p^t + p_{\text{RT}}^t \leq p_{\max}, t = 1, \dots, T, \quad (4c)$$

where  $p^1, \dots, p^T$  is the power allocated to the user during all the  $T$  time slots. (4b) is the transmission rate constraint of pre-downloading  $B$  bits within  $T$  time slots, (4c) is the power constraint of the SBSs.

The optimal solution of problem (4) satisfies the following multi-level water-filling structure [10]

$$p^{t*} = \begin{cases} \left( \frac{W^t}{W_{\max}} \nu^* - \frac{1}{g^t} \right)_0^{p_{\max} - p_{\text{RT}}^t}, & t \in \mathcal{T}_{\text{oc}} \cup \mathcal{N}^* \\ 0, & t \in \mathcal{T}_{\text{id}} - \mathcal{N}^* \end{cases}, \quad (5)$$

where  $\nu^*$  is the water-filling level,  $\mathcal{T}_{\text{oc}} = \{t | p_{\text{RT}}^t > 0\}$  or  $\mathcal{T}_{\text{id}} = \{t | p_{\text{RT}}^t = 0\}$  is the index set of the time slots that have or not have the RT traffic, and  $\mathcal{N}^* = \{t | g^t \geq g_{\text{th}}^*, t \in \mathcal{T}_{\text{id}}\}$  is index set of the scheduled  $N^*$  time slots without RT traffic for pre-downloading, which is determined by a threshold  $g_{\text{th}}^*$ , and the function  $(x)_0^a$  means  $\max\{\min\{a, x\}, 0\}$ . Since  $p^{t*} \geq 0$ ,  $t \in \mathcal{N}^*$  holds for any  $g^t \geq g_{\text{th}}^*$ , we have,

$$\nu^* - \frac{1}{g_{\text{th}}^*} \geq 0. \quad (6)$$

In practice, due to the user mobility and random arrival of the RT service, the information of  $g^t, W^t$  and  $p_{\text{RT}}^t$  in all the  $T$  time slots is hard to know, especially the small scale fading gains and residual resources available for pre-downloading at

each SBS. Fortunately, the information in all time slots is only needed when computing the water-filling level  $\nu^*$  and the threshold  $g_{\text{th}}^*$  [10]. If the MBS can estimate  $\nu^*$  and  $g_{\text{th}}^*$  from the following context information and then inform the SBS closest to the user, the SBS can allocate the power using (5) in the  $t$ th time slot only with the knowledge of  $g^t$ ,  $W^t$  and  $p_{\text{RT}}^t$  in this time slot.

- *Network level context information:*

From the statistics of traffic loads in the past, the resource utilization status of a SBS can be estimated. The statue can be modelled as a probability that  $(1 - \frac{1}{L}) \cdot 100\%$  bandwidth is occupied by RT traffic in each frame, i.e.,  $P_l^j \triangleq \Pr(W^t = \frac{1}{L} W_{\text{max}}), t = 1 + (j-1)T_s, \dots, jT_s, j = 1, \dots, T_f, l \in \{0, 1, \dots, L\}$ , where  $L$  is the maximal number of RT requests that a SBS can support in one time slot. The reserved transmit power for RT service is assumed in proportion to the occupied bandwidth, i.e.,  $p_{\text{RT}}^t = (1 - \frac{W^t}{W_{\text{max}}})p_{\text{max}}$ . We assume that  $P_l^j, l = 0, \dots, L$  are known at the SBS, which is reported to the MBS as the network level context information. Note that only  $P_L^j$  (i.e., the probability that a SBS is not occupied by RT traffic) is assumed known in [10].

- *User level context information:*

From the predicted trajectory of the user, the large scale fading gains can be obtained from the radio map as  $\alpha^1, \dots, \alpha^{T_f}$  during the  $T_f$  frames [12]. Consider  $g^t = \alpha^{\lceil \frac{t}{T_s} \rceil} \|\mathbf{h}^t\|^2 / (N_0 W^t)$  and denote  $\tilde{g}^t \triangleq \alpha^{\lceil \frac{t}{T_s} \rceil} \|\mathbf{h}^t\|^2 / (N_0 W_{\text{max}})$ , then the equivalent channel gain can be rewritten as  $g^t = \frac{W_{\text{max}}}{W^t} \tilde{g}^t$ . For Rayleigh fading channel,  $\tilde{g}^t$  in the  $j$ th frame follows Gamma distribution, whose probability density function (pdf) is

$$f^j(g) = \frac{1}{\Gamma(N_t)} \left( \frac{N_0 W_{\text{max}}}{\alpha^j} g \right)^{N_t-1} \exp \left( -\frac{N_0 W_{\text{max}}}{\alpha^j} g \right), \quad (7)$$

which is assumed known at the SBS, and is reported to the MBS as the user level context information.

To estimate the water-filling level and threshold with the two levels of context information, we transform problem (4) into another problem by using the relation of  $p^{t*}, t = 1, \dots, T$  with  $\nu^*$  and  $g_{\text{th}}^*$  in (5), and consider the case where the small-scale fading is ergodic in each frame.

In what follows, we separately transform the objective function and constraints of problem (4).

By substituting (5) into (4a), we can obtain the following proposition.

*Proposition 1:* When  $\nu \leq p_{\text{max}}$  and the small scale fading and available bandwidth are ergodic in each frame, the objective function in (4a) becomes

$$\begin{aligned} & \frac{1}{\xi} \frac{1}{T_f} \sum_{j=1}^{T_f} \sum_{l=1}^{L-1} P_l^j \frac{1}{L} \int_{\frac{1}{\nu}}^{\infty} \left( \nu - \frac{1}{g} \right) f^j(g) dg + \\ & \frac{1}{\xi} \frac{1}{T_f} \sum_{j=1}^{T_f} P_L^j \int_{g_{\text{th}}}^{\infty} \left( \nu - \frac{1}{g} \right) f^j(g) dg + \\ & \frac{p_{\text{act}} - p_{\text{sle}}}{T_f} \sum_{j=1}^{T_f} P_L^j \int_{g_{\text{th}}}^{\infty} f^j(g) dg. \end{aligned} \quad (8)$$

Using similar derivations, we can obtain Proposition 2.

*Proposition 2:* When  $\nu \leq p_{\text{max}}$  and the small scale fading and available bandwidth are ergodic in each frame, constraint (4b) becomes

$$\begin{aligned} & \frac{1}{T_f} \sum_{j=1}^{T_f} \sum_{l=1}^{L-1} P_l^j \frac{1}{L} W_{\text{max}} \int_{\frac{1}{\nu}}^{\infty} \ln(\nu g) f^j(g) dg + \\ & \frac{1}{T_f} \sum_{j=1}^{T_f} P_L^j W_{\text{max}} \int_{g_{\text{th}}}^{\infty} \ln(\nu g) f^j(g) dg = \frac{B \ln 2}{T \Delta_t}. \end{aligned} \quad (9)$$

The maximal and minimal power constraints in (4c) are guaranteed by the function  $(\cdot)_0^a$  in (5). Then, the water-filling level and threshold can be estimated from the following optimization problem,

$$\begin{aligned} \mathbf{P1} : & \min_{\nu, g_{\text{th}}} \quad (8) \\ & \text{s.t.} \quad (9), (6) \end{aligned}$$

Denote the solution of problem **P1** as  $\hat{\nu}^*$  and  $\hat{g}_{\text{th}}^*$ . Problem **P1** is equivalent to problem (4) when channel and available bandwidth are ergodic and when the solution of (4) satisfies  $\nu^* \leq p_{\text{max}}$ ,<sup>1</sup> in the sense that the optimal solutions of  $\nu$  and  $g_{\text{th}}$  obtained from the two problems are identical.

In **P1**, the objective function and constraints are differentiable with respect to the variables  $\nu$  and  $g_{\text{th}}$  since they are the integration of continuous functions. Hence, the optimal solution must satisfy the *Karush-Kuhn-Tucker* (KKT) conditions [13]. From the KKT conditions of problem **P1**, we can prove Proposition 3. The proof is omitted due the space limitation.

*Proposition 3:* The optimal water-filling level  $\hat{\nu}^*$  and threshold  $\hat{g}_{\text{th}}^*$  satisfy the following equation,

$$\left( \hat{\nu}^* - \frac{1}{\hat{g}_{\text{th}}^*} \right) + \xi(p_{\text{act}} - p_{\text{sle}}) - \hat{\nu}^* \ln(\hat{\nu}^* \hat{g}_{\text{th}}^*) = 0 \quad (11)$$

and  $\hat{\nu}^* \hat{g}_{\text{th}}^* > 1$ .

Further considering that the optimal solution should satisfy the constraint in (9),  $\hat{\nu}^*$  and  $\hat{g}_{\text{th}}^*$  satisfy two equalities in (9) and (11). By taking the derivation of (11) with respect to  $\hat{g}_{\text{th}}^*$  and considering  $\hat{\nu}^* \hat{g}_{\text{th}}^* > 1$  in Proposition 3, we have  $\frac{\partial \hat{\nu}^*}{\partial \hat{g}_{\text{th}}^*} = \frac{1}{\hat{g}_{\text{th}}^*} \left( \frac{1}{\hat{g}_{\text{th}}^*} - \hat{\nu}^* \right) / \ln(\hat{\nu}^* \hat{g}_{\text{th}}^*) < 0$ , i.e.,  $\hat{\nu}^*$  is a monotonic decreasing function of  $\hat{g}_{\text{th}}^*$ . Using the similar way, we can show that the left hand side of (9) is a monotonic decreasing function of  $\hat{g}_{\text{th}}^*$ . This suggests that we can find the global optimal solution of problem **P1** by a two-tier bisection searching algorithm. In the inner tier, we find  $\hat{\nu}^*$  with given  $\hat{g}_{\text{th}}^*$  by bisection searching from (11). In the outer tier, we find  $\hat{g}_{\text{th}}^*$  by bisection searching from (9).

The single-user pre-downloading strategy during the off-peak time can be summarized as the following two steps.

- 1) The MBS estimates  $\hat{\nu}^*$  and  $\hat{g}_{\text{th}}^*$  using the two-tier bisection searching algorithm based on the context information. Since the information is in long term, the role of such an estimation is to make the resource planning.
- 2) In the  $t$ th time slot, with the instantaneous information of  $g^t, W_{\text{RT}}^t, p_{\text{RT}}^t$ , the SBS optimizes the transmit power based on (5) with the estimated water-filling level and threshold, and then  $s^t$  can be obtained. When  $p^{t*} = 0$ , the user will not be scheduled by the SBS.

## B. Multi-user Pre-downloading

When targeting at minimizing the total energy consumption for pre-downloading files to multiple users, scheduling and

<sup>1</sup>The solution of problem (4) satisfies this condition easily since the duration  $T \Delta_t$  for pre-downloading  $B$  bits is long.

power allocation need to be jointly designed, whose optimal solution is hard to find due to the coupling among users. In the sequel, we propose a heuristic strategy, where the user scheduling and power allocation are separately designed.

Inspired by the single-user strategy, the power allocation to the  $T$  time slots for each MS is designed by using context information. Since there may exist multiple MSs in each cell, the resource utilization status for each MS should also reflect the resources occupied by pushing contents for other MSs. The pre-downloading strategy for multiple users can be implemented as follows.

- 1) *Power Allocation*: After gathering the predicted trajectories of all MSs in the macro cell, the MBS is aware of the set of users located in the  $i$ th small cell in the  $j$ th frame (denoted as  $\mathcal{K}_i^j$  with cardinality  $K_i^j$ ). By assuming that BS <sub>$i$</sub>  selects one of the users in  $\mathcal{K}_i^j$  with same probability, the time resources (i.e., time slots) available for each user in  $\mathcal{K}_i^j$  will reduce  $K_i^j$  times. Therefore, for MS <sub>$k$</sub> , the resource utilization status of its closest SBS (say, BS <sub>$i$</sub> ) can be modelled as,

$$\hat{P}_{k,l}^j = \frac{P_{k,l}^j}{K_i^j}, \quad l = 0, \dots, L, \quad k \in \mathcal{K}_i^j \quad (12)$$

with which the MBSs can estimate the water-filling level and threshold for MS <sub>$k$</sub>  as  $\hat{\nu}_k^*$  and  $\hat{g}_{th,k}^*$  from problem **P1**. Then, in the  $t$ th time slot, which is in  $\lceil \frac{t}{T_s} \rceil$ th frame, BS <sub>$i$</sub>  can optimize the power to the users in  $\mathcal{K}_i^{\lceil \frac{t}{T_s} \rceil}$  according to their own water-filling level  $\hat{\nu}_k^*$  and threshold  $\hat{g}_{th,k}^*$ , as well as their own instantaneous channel  $g_k^t$  and available resources  $W_{\max} - W_{i,RT}^t$  and  $p_{\max} - p_{i,RT}^t$  by using (5).

- 2) *Multi-user Scheduler*: In the  $t$ th time slot, according to the estimated water-filling levels and thresholds, the set of MSs in the  $i$ th cell who are allocated with non-zero powers is defined as a conflict set (denoted as  $\hat{\mathcal{K}}_i^t$  with cardinality  $\hat{K}_i^t$ ). To ensure the fairness among the MSs in using the residual resources in each SBS, random scheduler is employed, i.e., the MSs in  $\hat{\mathcal{K}}_i^t$  is scheduled by BS <sub>$i$</sub>  with the same probability of  $\frac{1}{\hat{K}_i^t}$ .

#### IV. SIMULATION RESULTS

In this section, we evaluate the network throughput and energy consumption of the proposed pre-downloading strategy, where the energy consumed both by content placement and by content delivery (due to cache miss) is taken into account. For comparison, we also simulate a broadcasting strategy and a traditional transmission strategy without pre-caching.

We consider  $M = 19$  small cells each with radius  $D = 50$  m hexagonally placed in a macro cell with radius 250 m. The path-loss model is  $30.5 + 36.7 \log_{10}(d)$ , where  $d$  is the distance between BS and MS in meter [14], and  $\sigma^2 = -165 + 10 \lg(W_{\max}) = -95$  dBm. The small scale channel is subject to Rayleigh block fading. The bandwidth is  $W_{\max} = 10$  MHz, and the maximal transmit power of SBS is  $p_{\max} = 0.2$  W. The circuit power consumptions in active and sleep modes are  $p_{\text{act}} = 3$  W and  $p_{\text{sle}} = 1$  W, respectively. The power amplifier efficiency is  $\xi = 8\%$  [11]. The duration of each time slot is  $\Delta_t = 10$  ms, and each frame contains  $T_s = 100$  time

slots. Each SBS serves RT service arrived with average rate of  $\lambda_{RT} = 0.2$  requests per time slot, and the service time follows exponential distribution with average two time slots.  $20\%W_{\max}$  and  $20\%p_{\max}$  are reserved for each RT request.

The file size is set as  $F = 30$  MBytes [1]. The file catalog  $\mathcal{N}_f$  that all the users in the macro cell may request contains  $N_f = 10000$  files, where the files are indexed according to popularity. The user requests follow Zipf distribution with parameter  $\beta_f \in [0, 1]$ , where  $\beta_f$  reflects the “peakiness” of the common interests of the users. To model the uncertainty of the prediction for each user’s personal interest, the number of files that each user may request is set as  $N_s = 100$ , which is a subset of the file catalog  $\mathcal{N}_f$  [9], and the subsets for different users may differ. Moreover, the request of each user follows Zipf distribution with parameter  $\beta_s \in [0, 1]$ .

Each SBS pushes files to 10 MSs in its cell. Each MS moves along a line with random direction and constant speed of 1 m/s in a small cell, with minimal distance of 5 m to 40 m from the SBS. During the peak time, the content delivery request arrives at the SBS with average rate  $\lambda_{CD} = 1.2$  Mbps per MS. Since simulating the whole day takes long time, we choose 120 s from the off-peak time to calculate the energy consumed by pushing and 60 s from the peak time to calculate the energy consumed by delivering the requested files. The results are obtained from 1000 Monte-Carlo trails, where the trajectory of each MS stays the same and the small-scale channel and RT requests vary in every trail. Unless otherwise specified, this simulation setup is used for all results.

The strategies to be compared are detailed as follows.

- Pre-caching with broadcast (with legend “Broadcast”): With the predicted common interests for all users in a macro cell, the MBS caches the most popular  $N_b^c = 10$  files from  $\mathcal{N}_f$  to each user by broadcast once a day during the off-peak time. Since the energy consumed by broadcasting is negligible, we ignore its energy consumption.
- Pushing with unicast (with legend “Unicast”): With the predicted personal interest of each user, the proposed strategy pushes the most possibly requested  $N_u^c = 10$  files from the  $N_s$  files to each user by unicast during the off-peak time. Since the energy consumed by resource planning is negligible, we ignore the energy consumption at the MBS.
- Baseline (with legend “Baseline”): The files are transmitted by the SBSs after they are requested by the MSs.

We first show that the conditions in propositions 1 and 2 are easy to satisfy. Fig. 2 shows the cumulative distribution function (CDF) of the water-filling level  $\nu^*$  and threshold  $g_{th}^*$  obtained from problem (4) and  $\hat{\nu}^*$ ,  $\hat{g}_{th}^*$  obtained from problem **P1**, where 100 MBits of content needs to pre-download to a user during 120 s. We can see that  $\nu^* < p_{\max}$  is easy to hold, and  $|\nu^* - \hat{\nu}^*|/\hat{\nu}^* < 5\%$  and  $|g_{th}^* - \hat{g}_{th}^*|/\hat{g}_{th}^* < 15\%$  even when  $T_s = 100$ , which is far from infinity.

In Fig. 3, we show the network throughput versus the average arrival rate of the content delivery traffic and  $\beta_f = \beta_s$ . We can see that the throughput gain of pushing files based on each user’s personal interest over pre-caching popular files to

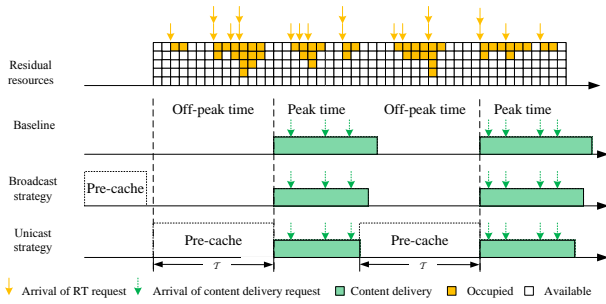


Fig. 1. Illustration of the simulated strategies.

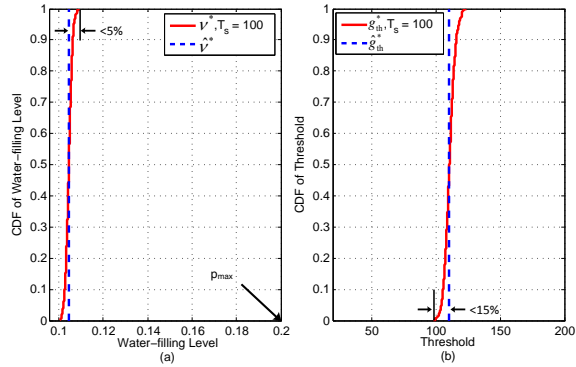


Fig. 2. (a) CDF of  $\hat{\nu}^*$  and  $\nu^*$ , (b) CDF of  $\hat{g}_{th}^*$  and  $g_{th}^*$

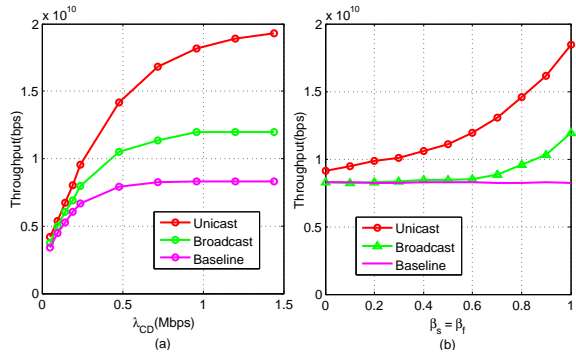


Fig. 3. Throughput vs. (a)  $\lambda_{CD}$ ,  $\beta_f = \beta_s = 1$  and (b)  $\beta_f = \beta_s$ ,  $\lambda_{CD} = 1.2$  Mbps

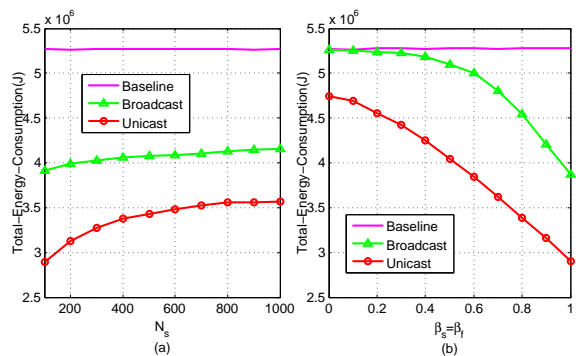


Fig. 4. Energy consumption vs. (a)  $N_s$ ,  $\beta_f = \beta_s = 1$ ,  $N_f = 100N_s$  and (b)  $\beta_f = \beta_s$ ,  $N_s = 100$

all users increases with  $\lambda_{CD}$ . As expected, the throughput gain from “Broadcast” and “Unicast” over the baseline increases

with  $\beta_f = \beta_s$ . When  $\beta_f$  is small, “Broadcast” achieves the same throughput as the “Baseline”, while “Unicast” can improve throughput in all cases.

In Fig. 4, we show the total energy consumed for content placement and delivery versus  $N_s$  and  $\beta_f = \beta_s$ . We can observe a surprising result that the proposed unicast strategy consumes less energy than broadcast. This is because more users can fetch their requested files in their own caches with pushing. Yet the energy-saving gain reduces with  $N_s$  since a large value of  $N_s$  indicates a high prediction uncertainty of each user’s interest, which leads to high cache miss rate and hence more energy for delivering the files. More energy can be saved by “Broadcast” and “Unicast” for larger  $\beta_f = \beta_s$  with respect to the “Baseline”, because the files cached at users have higher probability to be requested. When  $\beta_f$  is small, the “Broadcast” strategy consumes the same energy as the “Baseline”, while the proposed “Unicast” strategy can save energy in all cases.

## V. CONCLUSION

In this paper, we proposed a context-aware unicast pushing strategy based on the personal interest of each user, aimed at minimizing the extra energy consumed at the BSs for the file pre-downloading. Simulation results showed that the proposed strategy can improve network throughput and save energy compared to pre-caching with broadcast and traditional network without local caching when the files popularity is not uniform, even with the demand prediction uncertainty.

## REFERENCES

- [1] N. Golrezaei, A. F. Molisch, A. G. Dimakis, and G. Caire, “Femtocaching and device-to-device collaboration: A new architecture for wireless video distribution,” *IEEE Commun. Mag.*, vol. 51, no. 4, pp. 142–149, 2013.
- [2] H. Ahleghagh and S. Dey, “Video-aware scheduling and caching in the radio access network,” *IEEE/ACM Trans. Netw.*, vol. 22, no. 5, pp. 1444–1462, 2014.
- [3] D. Liu and C. Yang, “Will caching at base station improve energy efficiency of downlink transmission?” in *IEEE GlobalSIP*, 2014.
- [4] Y. Bao, X. Wang, S. Zhou, and Z. Niu, “An energy-efficient client pre-caching scheme with wireless multicast for video-on-demand services,” in *IEEE APCC*, 2012.
- [5] C. Yang, Z. Chen, Y. Yao, B. Xia, and H. Liu, “Energy efficiency in wireless cooperative caching networks,” in *IEEE ICC*, 2014.
- [6] K. Wang, Z. Chen, and H. Liu, “Push-based wireless converged networks for massive multimedia content delivery,” *IEEE Trans. Wireless Commun.*, vol. 13, no. 5, pp. 2894–2905, 2014.
- [7] B. Chen and C. Yang, “Performance gain of precaching at users in small cell networks,” in *IEEE PIMRC*, 2015.
- [8] B. D. Higgins, J. Flinn, T. J. Giulio, B. Noble, C. Peplin, and D. Watson, “Informed mobile prefetching,” in *ACM MobiSys*, 2012.
- [9] P. Lungaro, Z. Segall, and J. Zander, “Context-aware RRM for opportunistic content delivery in cellular networks,” in *IEEE CTRQ*, 2010.
- [10] C. Yao, C. Yang, and Z. Xiong, “Power-saving resource allocation by exploiting the context information,” in *IEEE PIMRC*, 2015.
- [11] G. Auer, O. Blume, V. Giannini, I. Gódor, et al., “D 2.3: Energy efficiency analysis of the reference systems, areas of improvements and target breakdown,” *EARTH*, Nov. 2010. [Online]. Available: <https://www.ict-earth.eu/publications/deliverables/deliverables.html>
- [12] H. Abou-Zeid and H. S. Hassanein, “Toward green media delivery: location-aware opportunities and approaches,” *IEEE Wireless Commun.*, vol. 21, no. 4, pp. 38–46, Aug. 2014.
- [13] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge University Press, 2004.
- [14] TR 36.814 V1.2.0, “Further Advancements for E-UTRA Physical Layer Aspects (Release 9),” *3GPP*, June 2009.